

Absorbing random-walk centrality: Theory and algorithms

Charalampos Mavroforakis
Dept. of Computer Science
Boston University
Boston, U.S.A.
cmav@cs.bu.edu

Michael Mathioudakis and Aristides Gionis
Helsinki Institute for Information Technology HIIT
Dept. of Computer Science, Aalto University
Helsinki, Finland
firstname.lastname@aalto.fi

Abstract—We study a new notion of graph centrality based on absorbing random walks. Given a graph $G = (V, E)$ and a set of query nodes $Q \subseteq V$, we aim to identify the k most central nodes in G with respect to Q . Specifically, we consider central nodes to be *absorbing* for random walks that start at the query nodes Q . The goal is to find the set of k central nodes that minimizes the expected length of a random walk until absorption. The proposed measure, which we call *k absorbing random-walk centrality*, favors *diverse* sets, as it is beneficial to place the k absorbing nodes in different parts of the graph so as to “intercept” random walks that start from different query nodes.

Although similar problem definitions have been considered in the literature, e.g., in information-retrieval settings where the goal is to diversify web-search results, in this paper we study the problem formally and prove some of its properties. We show that the problem is NP-hard, while the objective function is monotone and supermodular, implying that a greedy algorithm provides solutions with an approximation guarantee. On the other hand, the greedy algorithm involves expensive matrix operations that make it prohibitive to employ on large datasets. To confront this challenge, we develop more efficient algorithms based on spectral clustering and on personalized PageRank.

Keywords—graph mining; node centrality; random walks

I. INTRODUCTION

A fundamental problem in graph mining is to identify the most central nodes in a graph. Numerous centrality measures have been proposed, including degree centrality, closeness centrality [14], betweenness centrality [5], random-walk centrality [13], Katz centrality [9], and PageRank [4].

In the interest of robustness many centrality measures use random walks: while the shortest-path distance between two nodes can change dramatically by inserting or deleting a single edge, distances based on random walks account for multiple paths and offer a more global view of the connectivity between two nodes. In this spirit, the random-walk centrality of *one* node with respect to *all nodes* of the graph is defined as the expected time needed to come across this node in a random walk that starts in any other node of the graph [13].

In this paper, we consider a measure that generalizes random-walk centrality for a *set* of nodes C with respect to a set of *query nodes* Q . Our centrality measure is defined as the expected length of a random walk that starts from any node in Q until it reaches any node in C — at which point the random walk is “*absorbed*” by C . Moreover, to allow for

adjustable importance of query nodes in the centrality measure, we consider random walks *with restarts*, that occur with a fixed probability α at each step of the random walk. The resulting computational problem is to find a set of k nodes C that optimizes this measure with respect to nodes Q , which are provided as input. We call this measure *k absorbing random-walk centrality* and the corresponding optimization problem *k-ARW-CENTRALITY*.

To motivate the *k-ARW-CENTRALITY* problem, let us consider the scenario of searching the Web graph and summarizing the search results. In this scenario, nodes of the graph correspond to webpages, edges between nodes correspond to links between pages, and the set of query nodes Q consists of all nodes that match a user query, i.e., all webpages that satisfy a keyword search. Assuming that the size of Q is large, the goal is to find the k most central nodes with respect to Q , and present those to the user.

It is clear that ordering the nodes of the graph by their individual random-walk centrality scores and taking the top- k set does not solve the *k-ARW-CENTRALITY* problem, as these nodes may all be located in the same “neighborhood” of the graph, and thus, may not provide a good absorbing set for the query. On the other hand, as the goal is to minimize the expected absorption time for walks starting at Q , the optimal solution to the *k-ARW-CENTRALITY* problem will be a set of k , both *centrally-placed* and *diverse*, nodes.

This observation has motivated researchers in the information-retrieval field to consider random walks with absorbing states in order to diversify web-search results [18]. However, despite the fact that similar problem definitions and algorithms have been considered earlier, the *k-ARW-CENTRALITY* problem has not been formally studied and there has not been a theoretical analysis of its properties.

Our key results in this paper are the following: we show that the *k-ARW-CENTRALITY* problem is NP-hard, and we show that the k absorbing random-walk centrality measure is *monotone* and *supermodular*. The latter property allows us to quantify the approximation guarantee obtained by a natural greedy algorithm, which has also been considered by previous work [18]. Furthermore, a naïve implementation of the greedy algorithm requires many expensive matrix inversions, which make the algorithm particularly slow. Part of our contribu-

tion is to show how to make use of the Sherman-Morrison inversion formula to implement the greedy algorithm with only one matrix inversion and more efficient matrix \times vector multiplications.

Moreover, we explore the performance of faster, heuristic algorithms, aiming to identify methods that are faster than the greedy approach without significant loss in the quality of results. The heuristic algorithms we consider include the personalized PageRank algorithm [4], [10] as well as algorithms based on spectral clustering [17]. We find that, in practice, the personalized PageRank algorithm offers a very good trade-off between speed and quality.

The rest of the paper is organized as follows. In Section II, we overview previous work and discuss how it compares to this paper. We define our problem in Section III and provide basic background results on absorbing random walks in Section IV. Our main technical contributions are given in Sections IV and V, where we characterize the complexity of the problem, and provide the details of the greedy algorithm and the heuristics we explore. We evaluate the performance of algorithms in Section VII, over a range of real-world graphs, and Section VIII is a short conclusion. Proofs for some of the theorems shown in the paper are provided in the Appendix.

II. RELATED WORK

Many works in the literature explore ways to quantify the notion of node centrality on graphs [3]. Some of the most commonly-used measures include the following: (i) *degree centrality*, where the centrality of a node is simply quantified by its degree; (ii) *closeness centrality* [11], [14], defined as the average distance of a node from all other nodes on the graph; (iii) *betweenness centrality* [5], defined as the number of shortest paths between pairs of nodes in the graph that pass through a given node; (iv) *eigenvector centrality*, defined as the stationary probability that a Markov chain on the graph visits a given node, with Katz centrality [9] and PageRank [4] being two well-studied variants; and (v) *random-walk centrality* [13], defined as the expected first passage time of a random walk from a given node, when it starts from a random node of the graph. The measure we study in this paper generalizes the notion of *random-walk centrality* to a set of absorbing nodes.

Absorbing random walks have been used in previous work to select a *diverse* set of nodes from a graph. For example, an algorithm proposed by Zhu et al. [18] selects nodes in the following manner: (i) the first node is selected based on its PageRank value and is set as absorbing; (ii) the next node to be selected is the node that maximizes the expected first-passage time from the already selected absorbing nodes. Our problem definition differs considerably from the one considered in that work, as in our work the expected first-passage times are always computed from the set of query nodes that are provided in the input, and not from the nodes that participate in the solution so far. In this respect, the greedy method proposed by Zhu et al. is not associated with a crisp problem definition.

Another conceptually related line of work aims to select a diverse subset of query results, mainly within the context of document retrieval [1], [2], [16]. The goal, there, is to select k query results to optimize a function that quantifies the trade-off between relevance and diversity.

Our work is also remotely related to the problem studied by Leskovec et al. on *cost-effective outbreak detection* [12]. One of the problems discussed there is to select nodes in the network so that the detection time for a set of cascades is minimized. However, their work differs from ours on the fact that they consider as input a set of *cascades*, each one of finite size, while in our case the input consists of a set of query *nodes* and we consider a probabilistic model that generates random walk paths, of possibly infinite size.

III. PROBLEM DEFINITION

We are given a graph $G = (V, E)$ over a set of nodes V and set of undirected edges E . The number of nodes $|V|$ is denoted by n and the number of edges $|E|$ by m . The input also includes a subset of nodes $Q \subseteq V$, to which we refer as the *query nodes*. As a special case, the set of query nodes Q may be equal to the whole set of nodes, i.e., $Q = V$.

Our goal is to find a set C of k nodes that are *central* with respect to the query nodes Q . For some applications it makes sense to restrict the central nodes to be only among the query nodes, while in other cases, the central nodes may include any node in V . To model those different scenarios, we consider a set of candidate nodes D , and require that the k central nodes should belong in this candidate set, i.e., $C \subseteq D$. Some of the cases include $D = Q$, $D = V$, or $D = V \setminus Q$, but it could also be that D is defined in some other way that does not involve Q . In general, we assume that D is given as input.

The *centrality* of a set of nodes C with respect to query nodes Q is based on the notion of absorbing random-walks and their expected length. More specifically, let us consider a random walk on the nodes V of the graph, that proceeds at discrete steps: the walk starts from a node $q \in Q$ and, at each step moves to a different node, following edges in G , until it arrives at some node in C . The *starting* node q of the walk is chosen according to a probability distribution s . When the walk arrives at a node $c \in C$ for the first time, it terminates, and we say that the random walk is *absorbed* by that node c . In the interest of generality, and to allow for adjustable importance of query nodes in the centrality measure, we also allow the random walk to restart. Restarts occur with a probability α at each step of the random walk, where α is a parameter that is specified as input to the problem. When restarting, the walk proceeds to a query node selected randomly according to s . Intuitively, larger values of α favor nodes that are closer to nodes Q .

We are interested in the expected length (i.e., number of steps) of the walk that starts from a query node $q \in Q$ until it gets absorbed by some node in C , and we denote this expected length by $ac_Q^q(C)$. We then define the *absorbing random-walk*

centrality of a set of nodes C with respect to query nodes Q , by

$$\text{ac}_Q(C) = \sum_{q \in Q} \mathbf{s}(q) \text{ac}_q^q(C).$$

The problem we consider in this paper is the following.

Problem 1 (k -ARW-CENTRALITY) *We are given a graph $G = (V, E)$, a set of query nodes $Q \subseteq V$, a set of candidate nodes $D \subseteq V$, a starting probability distribution \mathbf{s} over V such that $\mathbf{s}(v) = 0$ if $v \in V \setminus Q$, a restart probability α , and an integer k . We ask to find a set of k nodes $C \subseteq D$ that minimizes $\text{ac}_Q(C)$, i.e., the expected length of a random walk that starts from Q and proceeds until it gets absorbed in some node in C .*

In cases where we have no reason to distinguish among the query nodes, we consider the uniform starting probability distribution $\mathbf{s}(q) = 1/|Q|$. In fact, for simplicity of exposition, hereinafter we focus on the case of uniform distribution. However, we note that all our definitions and techniques generalize naturally, not only to general starting probability distributions $\mathbf{s}(q)$, but also to *directed* and *weighted* graphs.

IV. ABSORBING RANDOM WALKS

In this section we review some relevant background on absorbing random walks. Specifically, we discuss how to calculate the objective function $\text{ac}_Q(C)$ for Problem 1.

Let \mathbf{P} be the transition matrix for a random walk, with $\mathbf{P}(i, j)$ expressing the probability that the random walk will move to node j given that it is currently at node i . Since random walks can only move to absorbing nodes C , but not away from them, we set $\mathbf{P}(c, c) = 1$ and $\mathbf{P}(c, j) = 0$, if $j \neq c$, for all absorbing nodes $c \in C$. The set $T = V \setminus C$ of non-absorbing nodes is called *transient*. If $N(i)$ are the neighbors of a node $i \in T$ and $d_i = |N(i)|$ its degree, the transition probabilities from node i to other nodes are

$$\mathbf{P}(i, j) = \begin{cases} \alpha \mathbf{s}(j) & \text{if } j \in Q \setminus N(i), \\ (1 - \alpha)/d_i + \alpha \mathbf{s}(j) & \text{if } j \in N(i). \end{cases} \quad (1)$$

Here, \mathbf{s} represents the starting probability vector. For example, for the uniform distribution over query nodes we have $\mathbf{s}(i) = 1/|Q|$ if $i \in Q$ and 0 otherwise. The transition matrix of the random walk can be written as follows

$$\mathbf{P} = \begin{pmatrix} \mathbf{P}_{TT} & \mathbf{P}_{TC} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}. \quad (2)$$

In the equation above, \mathbf{I} is an $(n - |T|) \times (n - |T|)$ identity matrix and $\mathbf{0}$ a matrix with all its entries equal to 0; \mathbf{P}_{TT} is the $|T| \times |T|$ sub-matrix of \mathbf{P} that contains the transition probabilities between transient nodes; and \mathbf{P}_{TC} is the $|T| \times |C|$ sub-matrix of \mathbf{P} that contains the transition probabilities from transient to absorbing nodes.

The probability of the walk being on node j at exactly ℓ steps having started at node i , is given by the (i, j) -entry of the matrix \mathbf{P}_{TT}^ℓ . Therefore, the expected total number of times

that the random walk visits node j having started from node i is given by the (i, j) -entry of the $|T| \times |T|$ matrix

$$\mathbf{F} = \sum_{\ell=0}^{\infty} \mathbf{P}_{TT}^\ell = (\mathbf{I} - \mathbf{P}_{TT})^{-1}, \quad (3)$$

which is known as the *fundamental matrix* of the absorbing random walk. Allowing the possibility to start the random walk at an absorbing node (and being absorbed immediately), we see that the expected length of a random walk that starts from node i and gets absorbed by the set C is given by the i -th element of the following $n \times 1$ vector

$$\mathbf{L} = \mathbf{L}_C = \begin{pmatrix} \mathbf{F} \\ \mathbf{0} \end{pmatrix} \mathbf{1}, \quad (4)$$

where $\mathbf{1}$ is an $T \times 1$ vector of all 1s. We write $\mathbf{L} = \mathbf{L}_C$ to emphasize the dependence on the set of absorbing nodes C .

The expected number of steps when starting from a node in Q and until being absorbed by some node in C is then obtained by summing over all query nodes, i.e.,

$$\text{ac}_Q(C) = \mathbf{s}^T \mathbf{L}_C. \quad (5)$$

A. Efficient computation of absorbing centrality

Equation (5) pinpoints the difficulty of the problem we consider: even computing the objective function $\text{ac}_Q(C)$ for a candidate solution C requires an expensive matrix inversion; $\mathbf{F} = (\mathbf{I} - \mathbf{P}_{TT})^{-1}$. Furthermore, searching for the optimal set C involves an exponential number of candidate sets, while evaluating each one of them requires a matrix inversion.

In practice, we find that we can compute $\text{ac}_Q(C)$ much faster approximately, as shown in Algorithm 1. The algorithm follows from the infinite-sum expansion of Equation (5).

$$\begin{aligned} \text{ac}_Q(C) &= \mathbf{s}^T \mathbf{L}_C = \mathbf{s}^T \begin{pmatrix} \mathbf{F} \\ \mathbf{0} \end{pmatrix} \mathbf{1} = \mathbf{s}^T \left(\sum_{\ell=0}^{\infty} \mathbf{P}_{TT}^\ell \right) \mathbf{1} \\ &= \mathbf{s}^T \sum_{\ell=0}^{\infty} \begin{pmatrix} \mathbf{P}_{TT}^\ell \\ \mathbf{0} \end{pmatrix} \mathbf{1} = \left(\sum_{\ell=0}^{\infty} \mathbf{s}^T \begin{pmatrix} \mathbf{P}_{TT}^\ell \\ \mathbf{0} \end{pmatrix} \right) \mathbf{1} \\ &= \left(\sum_{\ell=0}^{\infty} \mathbf{x}_\ell \right) \mathbf{1} = \sum_{\ell=0}^{\infty} \mathbf{x}_\ell \mathbf{1}, \end{aligned}$$

with

$$\mathbf{x}_0 = \mathbf{s}^T \quad \text{and} \quad \mathbf{x}_{\ell+1} = \mathbf{x}_\ell \begin{pmatrix} \mathbf{P}_{TT} \\ \mathbf{0} \end{pmatrix}. \quad (6)$$

Note that computing each vector \mathbf{x}_ℓ requires time $\mathcal{O}(n^2)$. Algorithm 1 terminates when the increase of the sum due to the latest term falls below a pre-defined threshold ϵ .

V. PROBLEM CHARACTERIZATION

We now study the k -ARW-CENTRALITY problem in more detail. In particular, we show that the function ac_Q is monotone and supermodular, a property that is used later to provide an approximation guarantee for the greedy algorithm. We also show that k -ARW-CENTRALITY is NP-hard.

Algorithm 1 ApproximateAC

Input: Transition matrix \mathbf{P}_{TT} , threshold ϵ , starting probabilities \mathbf{s}

Output: Absorbing centrality ac_Q

$\mathbf{x}_0 \leftarrow \mathbf{s}^T$

$\delta \leftarrow \mathbf{x}_0 \cdot \mathbf{1}$

$ac \leftarrow \delta$

$\ell \leftarrow 0$

while $\delta < \epsilon$ **do**
 $\mathbf{x}_{\ell+1} \leftarrow \mathbf{x}_\ell \begin{pmatrix} \mathbf{P}_{TT} \\ \mathbf{0} \end{pmatrix}$

$\delta \leftarrow \mathbf{x}_{\ell+1} \cdot \mathbf{1}$

$ac \leftarrow ac + \delta$

$\ell \leftarrow \ell + 1$

return ac

Recall that a function $f : 2^V \rightarrow \mathbb{R}$ over subsets of a ground set V is *submodular* if it has the *diminishing returns* property

$$f(Y \cup \{u\}) - f(Y) \leq f(X \cup \{u\}) - f(X), \quad (7)$$

for all $X \subseteq Y \subseteq V$ and $u \notin Y$. The function f is *supermodular* if $-f$ is submodular. Submodularity (and supermodularity) is a very useful property for designing algorithms. For instance, minimizing a submodular function is a polynomial-time solvable problem, while the maximization problem is typically amenable to approximation algorithms, the exact guarantee of which depends on other properties of the function and requirements of the problem, e.g., monotonicity, matroid constraints, etc.

Even though the objective function $ac_Q(C)$ is given in closed-form by Equation (5), to prove its properties we find it more convenient to work with its descriptive definition, namely, $ac_Q(C)$ being the expected length for a random walk starting at nodes of Q before being absorbed at nodes of C .

For the rest of this section we consider that the set of query nodes Q is fixed, and for simplicity we write $ac = ac_Q$.

Proposition 1 (Monotonicity) *For all $X \subseteq Y \subseteq V$ it is $ac(Y) \leq ac(X)$.*

The proposition states that absorption time decreases with more absorbing nodes. The proof is given in the Appendix.

Next we show that the absorbing random-walk centrality measure $ac(\cdot)$ is supermodular.

Proposition 2 (Supermodularity) *For all sets $X \subseteq Y \subseteq V$ and $u \notin Y$ it is*

$$ac(X) - ac(X \cup \{u\}) \geq ac(Y) - ac(Y \cup \{u\}). \quad (8)$$

Proof: Given an instantiation of a random walk, we define the following propositions for any pair of nodes $i, j \in V$, non-negative integer ℓ , and set of nodes Z :

$A_{i,j}^\ell(Z)$: The random walk started at node i and visited node j after exactly ℓ steps, without visiting any node in set Z .

$B_{i,j}^\ell(Z, u)$: The random walk started at node i and visited node j after exactly ℓ steps, having previously visited node u but without visiting any node in the set Z .

It is easy to see that the set of random walks for which $A_{i,j}^\ell(Z)$ is `true` can be partitioned into those that visited u within the first ℓ steps and those that did not. Therefore, the probability that proposition $A_{i,j}^\ell(Z)$ is `true` for any instantiation of a random walk generated by our model is equal to

$$\Pr [A_{i,j}^\ell(Z)] = \Pr [A_{i,j}^\ell(Z \cup \{u\})] + \Pr [B_{i,j}^\ell(Z, u)]. \quad (9)$$

Now, let $\mathbf{\Lambda}(Z)$ be the number of steps for a random walk to reach the nodes in Z . $\mathbf{\Lambda}(Z)$ is a random variable and its expected value over all random walks generated by our model is equal to $ac(Z)$. Note that the proposition $\mathbf{\Lambda}(Z) \geq \ell + 1$ is `true` for a given instantiation of a random walk only if there is a pair of nodes $q \in Q$ and $j \in V \setminus Z$, for which the proposition $A_{q,j}^\ell(Z)$ is `true`. Therefore,

$$\Pr [\mathbf{\Lambda}(Z) \geq \ell + 1] = \sum_{q \in Q} \sum_{j \in V \setminus Z} \Pr [A_{q,j}^\ell(Z)]. \quad (10)$$

From the above, it is easy to calculate $ac(Z)$ as

$$\begin{aligned} ac(Z) &= E[\mathbf{\Lambda}(Z)] \\ &= \sum_{\ell=0}^{\infty} \ell \Pr [\mathbf{\Lambda}(Z) = \ell] \\ &= \sum_{\ell=1}^{\infty} \Pr [\mathbf{\Lambda}(Z) \geq \ell] \\ &= \sum_{\ell=0}^{\infty} \Pr [\mathbf{\Lambda}(Z) \geq \ell + 1] \\ &= \sum_{\ell=0}^{\infty} \sum_{q \in Q} \sum_{j \in V \setminus Z} \Pr [A_{q,j}^\ell(Z)]. \end{aligned} \quad (11)$$

The final property we will need is the observation that, for $X \subseteq Y$, $B_{i,j}^\ell(Y, u)$ implies $B_{i,j}^\ell(X, u)$ and thus

$$\Pr [B_{i,j}^\ell(X, u)] \geq \Pr [B_{i,j}^\ell(Y, u)]. \quad (12)$$

By using Equation (11), the Inequality (8) can be rewritten as

$$\begin{aligned} &\sum_{\ell=0}^{\infty} \sum_{q \in Q} \sum_{j \in V \setminus X} \Pr [A_{q,j}^\ell(X)] - \\ &\quad \sum_{\ell=0}^{\infty} \sum_{q \in Q} \sum_{j \in V \setminus \{X \cup \{u\}\}} \Pr [A_{q,j}^\ell(X \cup \{u\})] \\ &\geq \sum_{\ell=0}^{\infty} \sum_{q \in Q} \sum_{j \in V \setminus Y} \Pr [A_{q,j}^\ell(Y)] - \\ &\quad \sum_{\ell=0}^{\infty} \sum_{q \in Q} \sum_{j \in V \setminus \{Y \cup \{u\}\}} \Pr [A_{q,j}^\ell(Y \cup \{u\})]. \end{aligned} \quad (13)$$

We only need to show that the inequality holds for an arbitrary value of ℓ and $q \in Q$, that is

$$\begin{aligned} & \sum_{j \in V \setminus X} \Pr [A_{q,j}^\ell(X)] - \sum_{j \in V \setminus \{X \cup \{u\}\}} \Pr [A_{q,j}^\ell(X \cup \{u\})] \geq \\ & \sum_{j \in V \setminus Y} \Pr [A_{q,j}^\ell(Y)] - \sum_{j \in V \setminus \{Y \cup \{u\}\}} \Pr [A_{q,j}^\ell(Y \cup \{u\})]. \end{aligned} \quad (14)$$

Notice that $\Pr [A_{i,u}^\ell(Y \cup \{u\})] = 0$, so we can rewrite the above inequality as

$$\begin{aligned} & \sum_{j \in V \setminus X} \Pr [A_{q,j}^\ell(X)] - \sum_{j \in V \setminus X} \Pr [A_{q,j}^\ell(X \cup \{u\})] \geq \\ & \sum_{j \in V \setminus Y} \Pr [A_{q,j}^\ell(Y)] - \sum_{j \in V \setminus Y} \Pr [A_{q,j}^\ell(Y \cup \{u\})]. \end{aligned} \quad (15)$$

To show the latter inequality we start from the left hand side and use Inequality (12). We have

$$\begin{aligned} & \sum_{j \in V \setminus X} \Pr [A_{i,j}^\ell(X)] - \sum_{j \in V \setminus X} \Pr [A_{i,j}^\ell(X \cup \{u\})] \\ & = \sum_{j \in V \setminus X} \Pr [B_{i,j}^\ell(X, u)] \\ & \geq \sum_{j \in V \setminus Y} \Pr [B_{i,j}^\ell(Y, u)] \\ & = \sum_{j \in V \setminus Y} \Pr [A_{i,j}^\ell(Y)] - \sum_{j \in V \setminus Y} \Pr [A_{i,j}^\ell(Y \cup \{u\})], \end{aligned}$$

which completes the proof. \blacksquare

Finally, we establish the hardness of k absorbing centrality, defined in Problem 1.

Theorem 1 *The k -ARW-CENTRALITY problem is NP-hard.*

Proof: We obtain a reduction from the VERTEXCOVER problem [6]. An instance of the VERTEXCOVER problem is specified by a graph $G = (V, E)$ and an integer k , and asks whether there exists a set of nodes $C \subseteq V$ such that $|C| \leq k$ and C is a vertex cover, (i.e., for every $(i, j) \in E$ it is $\{i, j\} \cap C \neq \emptyset$). Let $|V| = n$.

Given an instance of the VERTEXCOVER problem, we construct an instance of the decision version of k -ARW-CENTRALITY by taking the same graph $G = (V, E)$ with query nodes $Q = V$ and asking whether there is a set of absorbing nodes C such that $|C| \leq k$ and $\text{ac}_Q(C) \leq 1 - \frac{k}{n}$.

We will show that C is a solution for VERTEXCOVER if and only if $\text{ac}_Q(C) \leq 1 - \frac{k}{n}$.

Assuming first that C is a vertex cover. Consider a random walk starting uniformly at random from a node $v \in Q = V$. If $v \in C$ then the length of the walk will be 0, as the walk will be absorbed immediately. This happens with probability $|C|/|V| = k/n$. Otherwise, if $v \notin C$ the length of the walk will be 1, as the walk will be absorbed in the next step (since C is a vertex cover all the neighbors of v need to belong in C). This happens with the rest of the probability $1 - k/n$.

Thus, the expected length of the random walk is

$$\text{ac}_Q(C) = 0 \cdot \frac{k}{n} + 1 \cdot \left(1 - \frac{k}{n}\right) = 1 - \frac{k}{n} \quad (16)$$

Conversely, assume that C is not a vertex cover for G . Then, there should be an uncovered edge (u, v) . A random walk that starts in u and then goes to v (or starts in v and then goes to u) will have length at least 2, and this happens with probability at least $\frac{2}{n} \frac{1}{d_{\max}} \geq \frac{2}{n^2}$. Then, following a similar reasoning as in the previous case, we have

$$\begin{aligned} \text{ac}_Q(C) & = \sum_{k=0}^{\infty} k \Pr(\text{absorbed in exactly } k \text{ steps}) \\ & = \sum_{k=1}^{\infty} \Pr(\text{absorbed after at least } k \text{ steps}) \\ & \geq \left(1 - \frac{k}{n}\right) + \frac{2}{n^2} > 1 - \frac{k}{n}. \end{aligned} \quad (17)$$

VI. ALGORITHMS

This section presents algorithms to solve the k -ARW-CENTRALITY problem. In all cases, the set of query nodes $Q \subseteq V$ is given as input, along with a set of candidate nodes $D \subseteq V$ and the restart probability α .

A. Greedy approach

The first algorithm is a standard greedy algorithm, denoted Greedy, which exploits the supermodularity of the absorbing random-walk centrality measure. It starts with the result set C equal to the empty set, and iteratively adds a node from the set of candidate nodes D , until k nodes are added. In each iteration the node added in the set C is the one that brings the largest improvement to ac_Q .

As shown before, the objective function to be minimized, i.e., ac_Q , is supermodular and monotonically decreasing. The Greedy algorithm is not an approximation algorithm for this minimization problem. However, it can be shown to provide an approximation guarantee for maximizing the *absorbing centrality gain* measure, defined below.

Definition 1 (Absorbing centrality gain) *Given a graph G , a set of query nodes Q , and a set of candidate nodes D , the absorbing centrality gain of a set of nodes $C \subseteq D$ is defined as*

$$\text{acg}_Q(C) = m_Q - \text{ac}_Q(C),$$

where $m_Q = \min_{v \in D} \{\text{ac}_Q(\{v\})\}$.

Justification of the gain function. The reason to define the absorbing centrality gain is to turn our problem into a submodular-maximization problem so that we can apply standard approximation-theory results and show that the greedy algorithm provides a constant-factor approximation guarantee. The *shift* m_Q quantifies the absorbing centrality of the best single node in the candidate set. Thus, the value of $\text{acg}_Q(C)$ expresses how much we gain in expected random-walk length

when we use the set C as absorbing nodes compared to when we use the best single node. Our goal is to maximize this gain.

Observe that the gain function acg_Q is not non-negative everywhere. Take for example any node u such that $\text{ac}_Q(\{u\}) > m_Q$. Then, $\text{acg}_Q(\{u\}) < 0$. Note also that we could have obtained a non-negative gain function by defining gain with respect to the *worst* single node, instead of the best. In other words, the gain function $\text{acg}'_Q(C) = M_Q - \text{ac}_Q(C)$, with $M_Q = \max_{v \in D} \{\text{ac}_Q(\{v\})\}$, is non-negative everywhere.

Nevertheless, the reason we use the gain function acg_Q instead of acg'_Q is that acg'_Q takes much larger values than acg_Q , and thus, a multiplicative approximation guarantee on acg'_Q is a weaker result than a multiplicative approximation guarantee on acg_Q . On the other hand, our definition of acg_Q creates a technical difficulty with the approximation guarantee, that is defined for non-negative functions. Luckily, this difficulty can be overcome easily by noting that, due to the monotonicity of acg_Q , for any $k > 1$, the optimal solution of the function acg_Q , as well as the solution returned by Greedy, are both non-negative.

Approximation guarantee. The fact that the Greedy algorithm gives an approximation guarantee to the problem of maximizing absorbing centrality gain is a standard result from the theory of submodular functions.

Proposition 3 *The function acg_Q is monotonically increasing, and submodular.*

Proposition 4 *Let $k > 1$. For the problem of finding a set $C \subseteq D$ with $|C| \leq k$, such that $\text{acg}_Q(C)$ is maximized, the Greedy algorithm gives a $(1 - \frac{1}{e})$ -approximation guarantee.*

We now discuss the complexity of the Greedy algorithm. A naïve implementation requires computing the absorbing centrality $\text{ac}_Q(C)$ using Equation (5) for each set C that needs to be evaluated during the execution of the algorithm. However, applying Equation (5) involves a matrix inversion, which is a very expensive operation. Furthermore, the number of times that we need to evaluate $\text{ac}_Q(C)$ is $\mathcal{O}(k|D|)$, as for each iteration of the greedy we need to evaluate the improvement over the current set of each of the $\mathcal{O}(|D|)$ candidates. The number of candidates can be very large, e.g., $|D| = n$, yielding an $\mathcal{O}(kn^4)$ algorithm, which is prohibitively expensive.

We can show, however, that we can execute Greedy significantly more efficiently. Specifically, we can prove the following two propositions.

Proposition 5 *Let C_{i-1} be a set of $i - 1$ absorbing nodes, \mathbf{P}_{i-1} the corresponding transition matrix, and let $\mathbf{F}_{i-1} = (\mathbf{I} - \mathbf{P}_{i-1})^{-1}$. Let $C_i = C_{i-1} \cup \{u\}$. Given \mathbf{F}_{i-1} the value $\text{ac}_Q(C_i)$ can be computed in $\mathcal{O}(n^2)$.*

Proposition 6 *Let C be a set of absorbing nodes, \mathbf{P} the corresponding transition matrix, and $\mathbf{F} = (\mathbf{I} - \mathbf{P})^{-1}$. Let $C' = C - \{v\} \cup \{u\}$, $u, v \in C$. Given \mathbf{F} the value $\text{ac}_Q(C')$ can be computed in time $\mathcal{O}(n^2)$.*

The proofs of these two propositions can be found in the Ap-

Algorithm 2 Greedy

Input: graph G , query nodes Q , candidates D , $k \geq 1$

Output: a set of k nodes C

Compute $\text{ac}_Q(\{v\})$ for arbitrary $v \in D$

For each $u \in (D - \{v\})$, use Prop.6 to compute $\text{ac}_Q(u)$

Select $u_1 \in D$ s.t. $u_1 \leftarrow \arg \max_{u \in D} \text{ac}_Q(u)$

Initialize solution $C \leftarrow \{u_1\}$

for $i = 2..k$ **do**

 For each $u \in D$, use Prop.5 to compute $\text{ac}_Q(C \cup \{u\})$

 Select $u_i \in D$ s.t. $u_i \leftarrow \arg \max_{u_i \in (D - C)} \text{ac}_Q(C \cup \{u_i\})$

 Update solution $C \leftarrow C \cup \{u_i\}$

return C

pendix. Proposition 5 implies that in order to compute $\text{ac}_Q(C_i)$ for absorbing nodes C_i in $\mathcal{O}(n^2)$, it is enough to maintain the matrix \mathbf{F}_{i-1} , computed in the previous step of the greedy algorithm for absorbing nodes C_{i-1} . Proposition 6, on the other hand, implies that we can compute the absorbing centrality of each set of absorbing nodes of a fixed size i in $\mathcal{O}(n^2)$, given the matrix \mathbf{F} , which is computed for one arbitrary set of absorbing nodes C of size i . Combined, the two propositions above yield a greedy algorithm that runs in $\mathcal{O}(kn^3)$ and offers the approximation guarantee discussed above. We outline it as Algorithm 2.

Practical speed-up. We found that the following heuristic lets us speed-up Greedy even further, with no significant loss in the quality of results. To select the first node for the solution set C (see Algorithm 2), we calculate the *PageRank* values of all nodes in D and evaluate ac_Q only for the $t \ll k$ nodes with highest PageRank score, where t is a fixed parameter. In what follows, we will be using this heuristic version of Greedy, unless explicitly stated otherwise.

B. Efficient heuristics

Even though Greedy runs in polynomial time, it can be quite inefficient when employed on moderately sized datasets (more than some tens of thousands of nodes). We thus describe algorithms that we study as efficient heuristics for the problem. These algorithms do not offer guarantee for their performance.

Spectral methods have been used extensively for the problem of graph partitioning. Motivated by the wide applicability of this family of algorithms, here we explore three spectral algorithms: SpectralQ, SpectralC, and SpectralD. We start by a brief overview of the spectral method; a comprehensive presentation can be found in the tutorial by von Luxburg [17].

The main idea of spectral approaches is to project the original graph into a low-dimensional Euclidean space so that distances between nodes in the graph correspond to Euclidean distances between the corresponding projected points. A standard spectral embedding method, proposed by Shi and Malik [15], uses the “random-walk” *Laplacian* matrix $\mathbf{L}_G = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A}$ of a graph G , where \mathbf{A} is the adjacency matrix of the graph, and forms the matrix $\mathbf{U} = [u_2, \dots, u_{d+1}]$ whose columns are the eigenvectors of \mathbf{L}_G that correspond to

the smallest eigenvalues $\lambda_2 \leq \dots \leq \lambda_{d+1}$, with d being the target dimension of the projection. The spectral embedding is then defined by mapping the i -th node of the graph to a point in \mathbb{R}^d , which is the i -row of the matrix \mathbf{U} .

The algorithms we explore are adaptations of the spectral method. They all start by computing the spectral embedding $\phi : V \rightarrow \mathbb{R}^d$, as described above, and then, proceed as follows: SpectralQ performs k -means clustering on the embeddings of the *query nodes*, where k is the desired size of the result set. Subsequently, it selects *candidate nodes* that are close to the computed centroids. Specifically, if s_i is the size of the i -th cluster, then k_i candidate nodes are selected whose embedding is the nearest to the i -th centroid. The number k_i is selected so that $k_i \propto s_i$ and $\sum k_i = k$.

SpectralC is similar to SpectralQ, but it performs the k -means clustering on the embeddings of the *candidate nodes*, instead of the query nodes.

SpectralD performs k -means clustering on the embeddings of the *query nodes*, where k is the desired result-set size. Then, it selects the k candidate nodes whose embeddings minimize the sum of squared ℓ_2 -distances from the centroids, with no consideration of the relative sizes of the clusters.

Personalized Pagerank (PPR). This is the standard Pagerank [4] algorithm with a damping factor equal to the restart probability α of the random walk and personalization probabilities equal to the start probabilities $\mathbf{s}(q)$. Algorithm PPR returns the k nodes with highest PageRank values.

Degree and distance centrality. Finally, we consider the standard degree and distance centrality measures.

Degree returns the k highest-degree nodes. Note that this baseline is oblivious to the query nodes.

Distance returns the k nodes with highest distance centrality with respect to Q . The distance centrality of a node u is defined as $\text{dc}(u) = \left(\sum_{v \in Q} d(u, v) \right)^{-1}$.

VII. EXPERIMENTAL EVALUATION

A. Datasets

We evaluate the algorithms described in Section VI on two sets of real graphs: one set of small graphs that allows us to compare the performance of the fast heuristics against the greedy approach; and one set of larger graphs, to compare the performance of the heuristics against each other on datasets of larger scale. Note that the bottleneck of the computation lies in the evaluation of centrality. Even though the technique we describe in Section IV-A allows it to scale to datasets of tens of thousands of nodes on a single processor, it is still prohibitively expensive for massive graphs. Still, our experimentation allows us to discover the traits of the different algorithms and understand what performance to anticipate when they are employed on graphs of massive size.

The datasets are listed in Table I. Small graphs are obtained from Mark Newman’s repository¹, larger graphs from SNAP.²

TABLE I: Dataset statistics

Dataset	$ V $	$ E $
karate	34	78
dolphins	62	159
lesmis	77	254
adjnoun	112	425
football	115	613
kddCoauthors	2891	2891
livejournal	3645	4141
ca-GrQc	5242	14496
ca-HepTh	9877	25998
roadnet	10199	13932
oregon-1	11174	23409

For kddCoauthors, livejournal, and roadnet we use samples of the original datasets. In the interest of repeatability, our code and datasets are made publicly available.³

B. Evaluation Methodology

Each experiment in our evaluation framework is defined by a graph G , a set of query nodes Q , a set of candidate nodes D , and an algorithm to solve the problem. We evaluate all algorithms presented in Section VI. For the set of candidate nodes D , we consider two cases: it is equal to either the set of query nodes, i.e., $D = Q$, or the set of all nodes, i.e., $D = V$.

Query nodes Q are selected randomly, using the following process: First, we select a set S of s seed nodes, uniformly at random among all nodes. Then, we select a ball $B(v, r)$ of predetermined radius $r = 2$, around each seed $v \in S$.⁴ Finally, from all balls, we select a set of query nodes Q of predetermined size q , with $q = 10$ and $q = 20$, respectively, for the small and larger datasets. Selection is done uniformly at random.

Finally, the restart probability α is set to $\alpha = 0.15$ and the starting probabilities \mathbf{s} are uniform over Q .

C. Implementation

All algorithms are implemented in Python using the NetworkX package [8], and were run on an Intel Xeon 2.83GHz with 32GB RAM.

D. Results

Figure 1 shows the centrality scores achieved by different algorithms on the small graphs for varying k (note: lower is better). We present two settings: on the left, the candidates are all nodes ($D = V$), and on the right, the candidates are only the query nodes ($D = Q$). We observe that PPR tracks well the quality of solutions returned by Greedy, while Degree and Distance often come close to that. Spectral algorithms do not perform that well.

Figure 2 is similar to Figure 1, but results on the larger datasets are shown, not including Greedy. When all nodes are candidates, PPR typically has the best performance, followed by Distance, while Degree is unreliable. The spectral algorithms typically perform worse than PPR.

¹<http://www-personal.umich.edu/~7Eemjn/netdata/>

²<http://snap.stanford.edu/data/index.html>

³<https://github.com/harrymvr/absorbing-centrality>

⁴For the planar roadnet dataset we use $r = 3$.

When only query nodes are candidates, all algorithms demonstrate similar performance, which is most typically worse than the performance of PPR (the best performing algorithm) in the previous setting. Both observations can be explained by the fact that the selection is very restricted by the requirement $D = Q$, and there is not much flexibility for the best performing algorithms to produce a better solution.

In terms of running time on the larger graphs, Distance returns within a few minutes (with observed times between 15 seconds to 5 minutes) while Degree returns within seconds (all observed times were less than 1 minute). Finally, even though Greedy returns within 1-2 seconds for the small datasets, it does not scale well for the larger datasets (running time is orders of magnitude worse than the heuristics and not included in the experiments).

Based on the above, we conclude that PPR offers the best trade-off of quality versus running time for datasets of at least moderate size (more than 10k nodes).

VIII. CONCLUSIONS

In this paper, we have addressed the problem of finding central nodes in a graph with respect to a set of query nodes Q . Our measure is based on absorbing random walks: we seek to compute k nodes that minimize the expected number of steps that a random walk will need to reach at (and be “absorbed” by) when it starts from the query nodes. We have shown that the problem is NP-hard and described an $\mathcal{O}(kn^3)$ greedy algorithm to solve it approximately. Moreover, we experimented with heuristic algorithms to solve the problem on large graphs. Our results show that, in practice, personalized PageRank offers a good combination of quality and speed.

REFERENCES

- [1] R. Agrawal, S. Gollapudi, A. Halverson, and S. Jeong. Diversifying search results. In *Proceedings of the Second ACM International Conference on Web Search and Data Mining*, pages 5–14. ACM, 2009.
- [2] A. Angel and N. Koudas. *Efficient diversity-aware search*. ACM, June 2011.
- [3] P. Boldi and S. Vigna. Axioms for centrality. *Internet Mathematics*, 2014.
- [4] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30, 1998.
- [5] L. Freeman. A set of measures of centrality based upon betweenness. *Sociometry*, 40, 1977.
- [6] M. Garey and D. Johnson. *Computers and intractability; A guide to the theory of NP-completeness*. W. H. Freeman & Co., 1990.
- [7] G. H. Golub and C. F. Van Loan. *Matrix computations*, volume 3. JHU Press, 2012.
- [8] A. A. Hagberg, D. A. Schult, and P. J. Swart. Exploring network structure, dynamics, and function using NetworkX. In *SciPy*, 2008.
- [9] L. Katz. A New Status Index Derived from Sociometric Index. *Psychometrika*, 1953.
- [10] A. N. Langville and C. D. Meyer. A survey of eigenvector methods for web information retrieval. *SIAM review*, 47(1):135–161, 2005.
- [11] H. J. Leavitt. Some effects of certain communication patterns on group performance. *The Journal of Abnormal and Social Psychology*, 1951.
- [12] J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, and N. Glance. Cost-effective outbreak detection in networks. In *SIGKDD*. ACM, 2007.
- [13] J. D. Noh and H. Rieger. Random walks on complex networks. *Phys. Rev. Lett.*, 92, 2004.
- [14] G. Sabidussi. The centrality index of a graph. *Psychometrika*, 31, 1966.

- [15] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [16] M. R. Vieira, H. L. Razente, M. C. N. Barioni, M. Hadjieleftheriou, D. Srivastava, A. J. M. Traina, and V. J. Tsotras. On query result diversification. *2011 IEEE International Conference on Data Engineering*, pages 1163–1174, 2011.
- [17] U. Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [18] X. Zhu, A. Goldberg, J. Van Gael, and D. Andrzejewski. Improving diversity in web search results re-ranking using absorbing random walks. In *NAACL-HLT*, 2007.

APPENDIX

A. Proposition 1

Proposition (Monotonicity) For all $X \subseteq Y \subseteq V$ it is $ac(Y) \leq ac(X)$.

Proof: Write G_X for the input graph G where the set X are absorbing nodes. Define G_Y similarly. Let $Z = Y \setminus X$. Consider a path p in G_X drawn from the distribution induced by the random walks on G_X . Let $\Pr[p]$ be the probability of the path and $\ell(p)$ its length. Let $\mathcal{P}(X)$ and $\mathcal{P}(Y)$ be the set of paths on G_X and G_Y . Finally, let $\mathcal{P}(Z, X)$ be the set of paths on G_X that pass from Z , and $\mathcal{P}(Z, X)$ the set of paths on G_X that do *not* pass from Z . We have

$$\begin{aligned} ac(X) &= \sum_{p \in \mathcal{P}(X)} \Pr[p] \ell(p) \\ &= \sum_{p \in \mathcal{P}(Z, X)} \Pr[p] \ell(p) + \sum_{p \in \mathcal{P}(Z, X)} \Pr[p] \ell(p) \\ &\geq \sum_{p \in \mathcal{P}(Y)} \Pr[p] \ell(p) \\ &= ac(Y), \end{aligned}$$

where the inequality comes from the fact that a path in G_X passing from Z and being absorbed by X corresponds to a shorter path in G_Y being absorbed by Y . ■

B. Proposition 5

Proposition Let C_{i-1} be a set of $i - 1$ absorbing nodes, \mathbf{P}_{i-1} the corresponding transition matrix, and $\mathbf{F}_{i-1} = (\mathbf{I} - \mathbf{P}_{i-1})^{-1}$. Let $C_i = C_{i-1} \cup \{u\}$. Given \mathbf{F}_{i-1} , the centrality score $ac_Q(C_i)$ can be computed in time $\mathcal{O}(n^2)$.

The proof makes use of the following lemma.

Lemma 1 (Sherman-Morrison Formula [7]) Let \mathbf{M} be a square $n \times n$ invertible matrix and \mathbf{M}^{-1} its inverse. Moreover, let \mathbf{a} and \mathbf{b} be any two column vectors of size n . Then, the following equation holds

$$(\mathbf{M} + \mathbf{ab}^T)^{-1} = \mathbf{M}^{-1} - \mathbf{M}^{-1} \mathbf{ab}^T \mathbf{M}^{-1} / (1 + \mathbf{b}^T \mathbf{M}^{-1} \mathbf{a}).$$

Proof: Without loss of generality, let the set of absorbing nodes be $C_{i-1} = \{1, 2, \dots, i - 1\}$. As in Section VI, the expected number of steps before absorption is given by the formulas

$$ac_Q(C_{i-1}) = \mathbf{s}_Q^T \mathbf{F}_{i-1} \mathbf{1},$$

with $\mathbf{F}_{i-1} = \mathbf{A}_{i-1}^{-1}$ and $\mathbf{A}_{i-1} = \mathbf{I} - \mathbf{P}_{i-1}$.

We proceed to show how to increase the set of absorbing nodes by one and calculate the new absorption time by updating \mathbf{F}_{i-1} in $\mathcal{O}(n^2)$. Without loss of generality, suppose we add node i to the absorbing nodes C_{i-1} , so that

$$C_i = C_{i-1} \cup \{i\} = \{1, 2, \dots, i-1, i\}.$$

Let \mathbf{P}_i be the transition matrix over G with absorbing nodes C_i . Like before, the expected absorption time by nodes C_i is given by the formulas

$$\text{ac}_Q(C_i) = \mathbf{s}_Q^T \mathbf{F}_i \mathbf{1},$$

$$\text{with } \mathbf{F}_i = \mathbf{A}_i^{-1} \text{ and } \mathbf{A}_i = \mathbf{I} - \mathbf{P}_i.$$

Notice that

$$\begin{aligned} \mathbf{A}_i - \mathbf{A}_{i-1} &= (\mathbf{I} - \mathbf{P}_i) - (\mathbf{I} - \mathbf{P}_{i-1}) = \mathbf{P}_{i-1} - \mathbf{P}_i \\ &= \begin{bmatrix} \mathbf{0}_{(i-1) \times n} \\ p_{i,1} \dots p_{i,n} \\ \mathbf{0}_{(n-i) \times n} \end{bmatrix} = \mathbf{a} \mathbf{b}^T \end{aligned}$$

where $p_{i,j}$ denotes the transition probability from node i to node j in transition matrix \mathbf{P}_{i-1} , and the column-vectors \mathbf{a} and \mathbf{b} are defined as

$$\begin{aligned} \mathbf{a} &= \begin{bmatrix} \overbrace{0 \dots 0}^{i-1} & 1 & \overbrace{0 \dots 0}^{n-i} \end{bmatrix}, \text{ and} \\ \mathbf{b} &= [p_{i,1} \dots p_{i,n}]. \end{aligned}$$

By a direct application of Lemma 1, it is easy to see that we can compute \mathbf{F}_i from \mathbf{F}_{i-1} with the following formula, at a cost of $\mathcal{O}(n^2)$ operations.

$$\mathbf{F}_i = \mathbf{F}_{i-1} - (\mathbf{F}_{i-1} \mathbf{a})(\mathbf{b}^T \mathbf{F}_{i-1}) / (1 + \mathbf{b}^T (\mathbf{F}_{i-1} \mathbf{a}))$$

We have thus shown that, given \mathbf{F}_{i-1} , we can compute \mathbf{F}_i , and therefore $\text{ac}_Q(C_i)$ as well, in $\mathcal{O}(n^2)$. ■

C. Proposition 6

Proposition Let C be a set of absorbing nodes, \mathbf{P} the corresponding transition matrix, and $\mathbf{F} = (\mathbf{I} - \mathbf{P})^{-1}$. Let $C' = C - \{v\} \cup \{u\}$, for $u, v \in C$. Given \mathbf{F} , the centrality score $\text{ac}_Q(C')$ can be computed in time $\mathcal{O}(n^2)$.

Proof: The proof is similar to the proof of Proposition 5. Without loss of generality, let the two sets of absorbing nodes be

$$\begin{aligned} C &= \{1, 2, \dots, i-1, i\}, \text{ and} \\ C' &= \{1, 2, \dots, i-1, i+1\}. \end{aligned}$$

Let \mathbf{P}' be the transition matrix with absorbing nodes C' . The absorbing centrality for the two sets of absorbing nodes C and C' is expressed as a function of the following two matrices

$$\mathbf{F} = \mathbf{A}^{-1}, \text{ with } \mathbf{A} = \mathbf{I} - \mathbf{P}, \text{ and}$$

$$\mathbf{F}' = \mathbf{A}'^{-1}, \text{ with } \mathbf{A}' = (\mathbf{I} - \mathbf{P}').$$

Notice that

$$\begin{aligned} \mathbf{A}' - \mathbf{A} &= (\mathbf{I} - \mathbf{P}') - (\mathbf{I} - \mathbf{P}) = \mathbf{P} - \mathbf{P}' \\ &= \begin{bmatrix} \mathbf{0}_{(i-1) \times n} \\ -p_{i,1} \dots -p_{i,n} \\ p_{i+1,0} \dots p_{i+1,n} \\ \mathbf{0}_{(n-i-1) \times n} \end{bmatrix} = \mathbf{a}_2 \mathbf{b}_2^T - \mathbf{a}_1 \mathbf{b}_1^T \end{aligned}$$

where $p_{i,j}$ denotes the transition probability from node i to node j in a transition matrix \mathbf{P}_0 where neither node i or $i+1$ is absorbing, and the column-vectors $\mathbf{a}_1, \mathbf{b}_1, \mathbf{a}_2, \mathbf{b}_2$ are defined as

$$\begin{aligned} \mathbf{a}_1 &= \begin{bmatrix} \overbrace{0 \dots 0}^{i-1} & 1 & 0 & \overbrace{0 \dots 0}^{n-i-1} \end{bmatrix} \\ \mathbf{b}_1 &= [p_{i,1} \dots p_{i,n}] \\ \mathbf{a}_2 &= \begin{bmatrix} \overbrace{0 \dots 0}^{i-1} & 0 & 1 & \overbrace{0 \dots 0}^{n-i-1} \end{bmatrix} \\ \mathbf{b}_2 &= [p_{i+1,1} \dots p_{i+1,n}]. \end{aligned}$$

By an argument similar with the one we made in the proof of Proposition 5, we can compute \mathbf{F}' in the following two steps from \mathbf{F} , each costing $\mathcal{O}(n^2)$ operations for the provided parenthesization

$$\begin{aligned} \mathbf{Z} &= \mathbf{F} - (\mathbf{Z} \mathbf{a}_2)(\mathbf{b}_2^T \mathbf{Z}) / (1 + \mathbf{b}_2^T (\mathbf{Z} \mathbf{a}_2)), \\ \mathbf{F}' &= \mathbf{Z} + (\mathbf{F} \mathbf{a}_1)(\mathbf{b}_1^T \mathbf{F}) / (1 + \mathbf{b}_1^T (\mathbf{F} \mathbf{a}_1)). \end{aligned}$$

We have thus shown that, given \mathbf{F} , we can compute \mathbf{F}' , and therefore $\text{ac}_Q(C')$ as well, in time $\mathcal{O}(n^2)$. ■

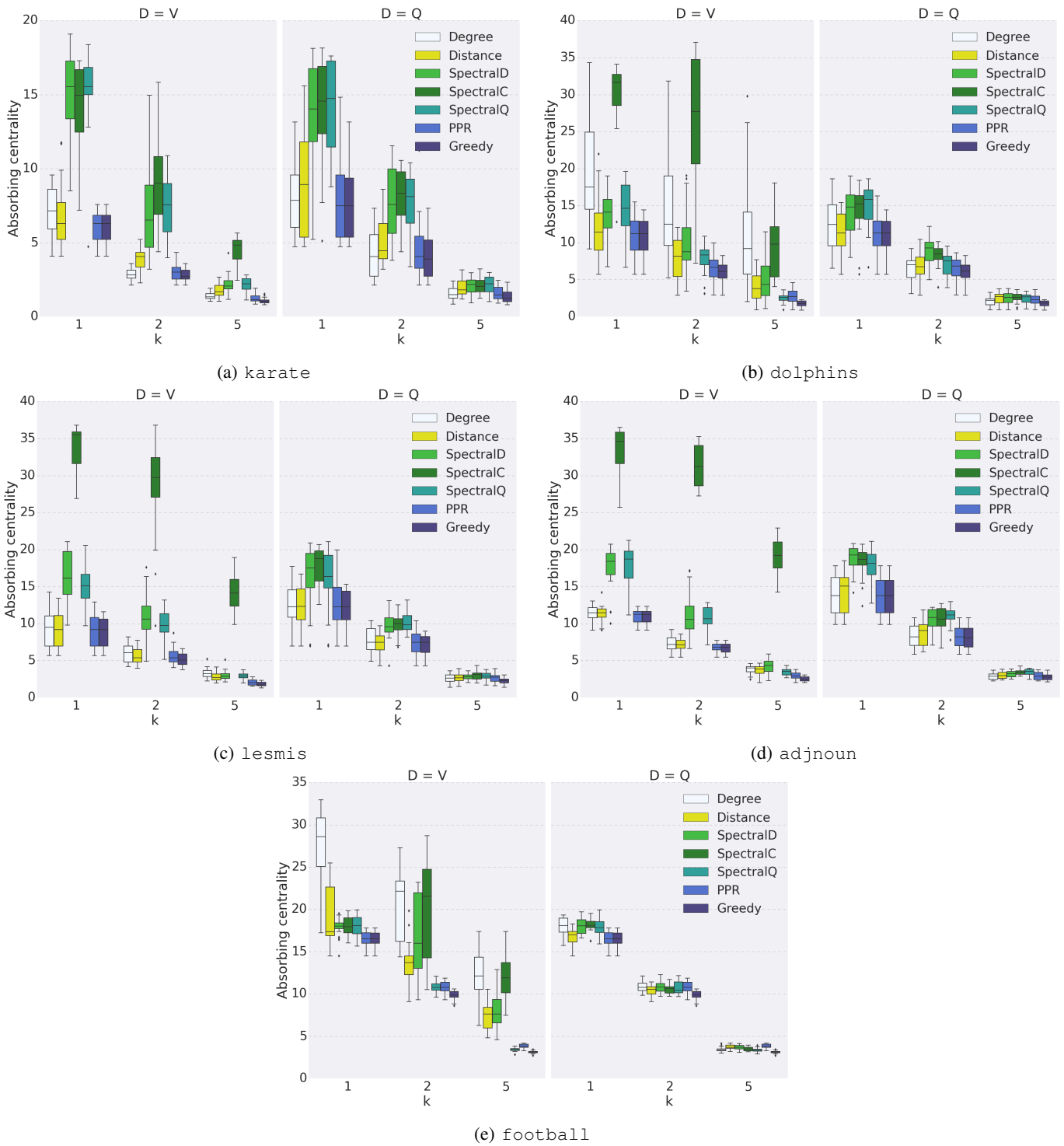


Fig. 1: Results on small datasets for varying k and $s = 2$.

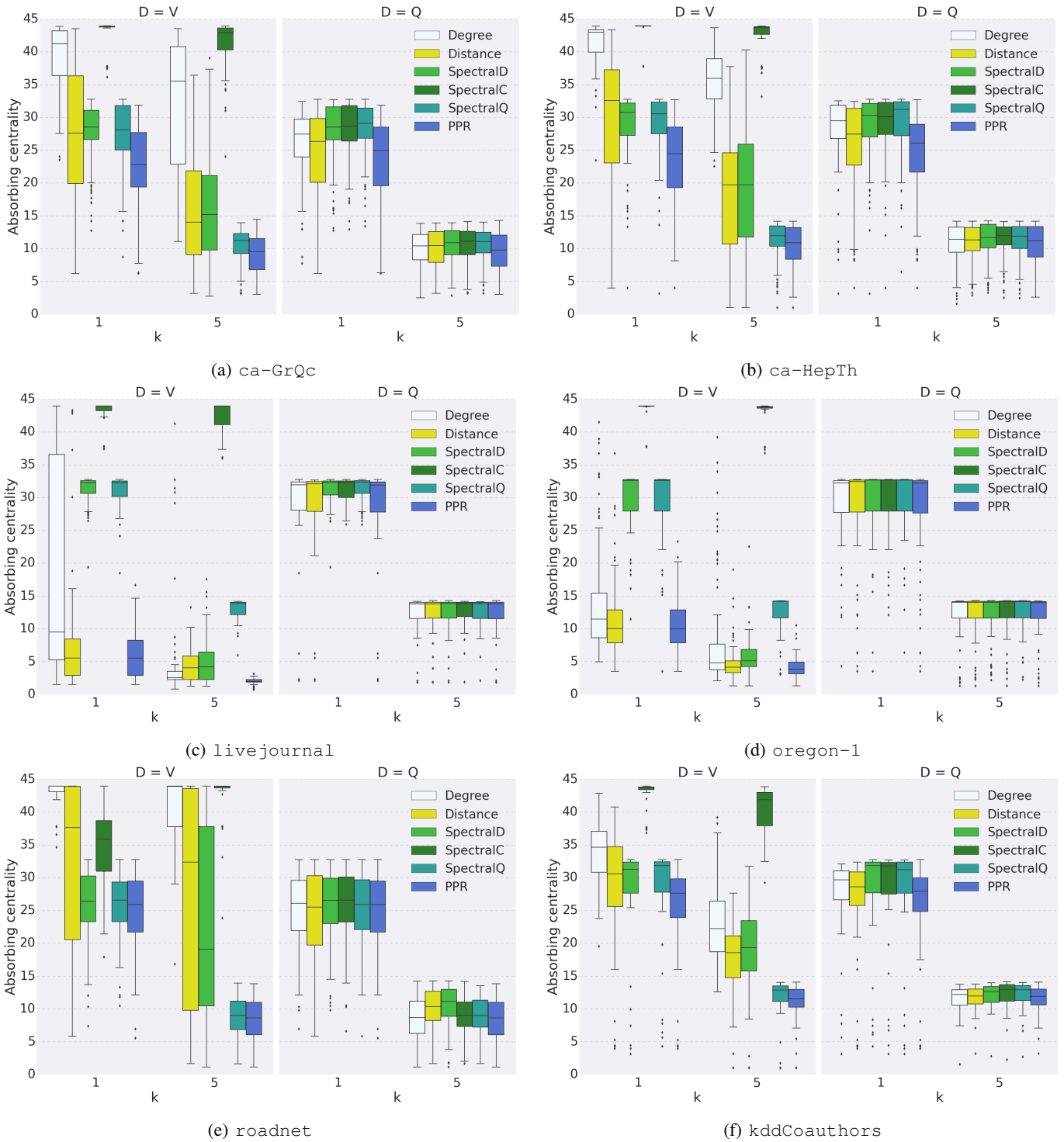


Fig. 2: Results on large datasets for varying k and $s = 5$.